

# Diclique clustering in a directed random graph

Mindaugas Bloznelis<sup>1</sup> and Lasse Leskelä<sup>2</sup>

<sup>1</sup> Vilnius University, Lithuania, [www.mif.vu.lt/~bloznelis/](http://www.mif.vu.lt/~bloznelis/)

<sup>2</sup> Aalto University, Finland, [math.aalto.fi/~lleskela/](http://math.aalto.fi/~lleskela/)

**Abstract.** We discuss a notion of clustering for directed graphs, which describes how likely two followers of a node are to follow a common target. The associated network motifs, called dicliques or bi-fans, have been found to be key structural components in various real-world networks. We introduce a two-mode statistical network model consisting of actors and auxiliary attributes, where an actor  $i$  decides to follow an actor  $j$  whenever  $i$  demands an attribute supplied by  $j$ . We show that the digraph admits nontrivial clustering properties of the aforementioned type, as well as power-law indegree and outdegree distributions.

**Keywords:** intersection graph, two-mode network, affiliation network, digraph, diclique, bi-fan, complex network

## 1 Introduction

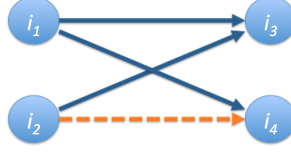
### 1.1 Clustering in directed networks

Many real networks display a tendency to cluster, that is, to form dense local neighborhoods in a globally sparse graph. In an undirected social network this may be phrased as: *your friends are likely to be friends*. This feature is typically quantified in terms of local and global clustering coefficients measuring how likely two neighbors of a node are neighbors [11,13,14,16]. In directed networks there are many ways to define the concept of clustering, for example by considering the thirteen different ways that a set of three nodes may form a weakly connected directed graph [5].

In this paper we discuss a new type of clustering concept which is motivated by directed online social networks, where a directed link  $i \rightarrow j$  means that an actor  $i$  follows actor  $j$ . In such networks a natural way to describe clustering is to say that *your followers are likely to follow common targets*. When the topology of the network is unknown and modeled as a random graph distributed according to a probability measure  $P$ , the above statement can be expressed as

$$P(i_2 \rightarrow i_4 \mid i_1 \rightarrow i_2, i_1 \rightarrow i_3, i_2 \rightarrow i_3) > P(i_2 \rightarrow i_4), \quad (1)$$

where 'you' corresponds to actor  $i_3$ . Interestingly, the conditional probability on the left can stay bounded away from zero even for sparse random digraphs [9]. The associated subgraph (Fig. 1) is called a *diclique*. Earlier experimental studies



**Fig. 1.** Forming a diclique by adding a link  $i_2 \rightarrow i_4$ .

have observed that dicliques (a.k.a. bi-fans) constitute a key structural motif in gene regulation networks [10], citation networks, and several types of online social networks [17].

Motivated by the above discussion, we define a global *diclique clustering coefficient* of a finite directed graph  $D$  with an adjacency matrix  $(D_{ij})$  by

$$\mathcal{C}_{\text{di}}(D) = \frac{\sum_{(i_1, i_2, i_3, i_4)} D_{i_1, i_3} D_{i_1, i_4} D_{i_2, i_3} D_{i_2, i_4}}{\sum_{(i_1, i_2, i_3, i_4)} D_{i_1, i_3} D_{i_1, i_4} D_{i_2, i_3}}, \quad (2)$$

where the sums are computed over all ordered quadruples of distinct nodes. It provides an empirical counterpart to the conditional probability (1) in the sense that the ratio in (2) defines the conditional probability

$$P_D(I_2 \rightarrow I_4 \mid I_1 \rightarrow I_3, I_1 \rightarrow I_4, I_2 \rightarrow I_3), \quad (3)$$

where  $P_D$  refers to the distribution of the random quadruple  $(I_1, I_2, I_3, I_4)$  sampled uniformly at random among all ordered quadruples of distinct nodes in  $D$ .

To quantify diclique clustering among the followers of a selected actor  $i$ , we may define a local diclique clustering coefficient by

$$\mathcal{C}_{\text{di}}(D, i) = \frac{\sum_{(i_1, i_2, i_4)} D_{i_1, i} D_{i_1, i_4} D_{i_2, i} D_{i_2, i_4}}{\sum_{(i_1, i_2, i_4)} D_{i_1, i} D_{i_1, i_4} D_{i_2, i}}, \quad (4)$$

where the sums are computed over all ordered triples of distinct nodes excluding  $i$ . We remark that  $\mathcal{C}_{\text{di}}(D, i) = P_D(I_2 \rightarrow I_4 \mid I_1 \rightarrow I_3, I_1 \rightarrow I_4, I_2 \rightarrow I_3, I_3 = i)$ .

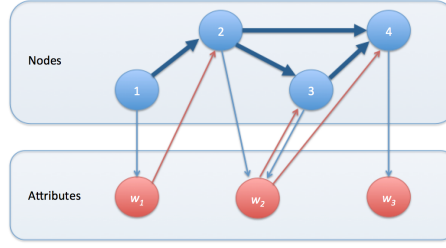
*Remark 1.* By replacing  $\rightarrow$  by  $\leftrightarrow$  in (3), we see that the analogue of the above notion for undirected graphs corresponds to predicting how likely the endpoints of the 3-path  $I_2 \leftrightarrow I_3 \leftrightarrow I_1 \leftrightarrow I_4$  are linked together.

## 1.2 A directed random graph model

Our goal is to define a parsimonious yet powerful statistical model of a directed social network which displays diclique clustering properties as discussed in the previous section. Clustering properties in many social networks, such as movie actor networks or scientific collaboration networks, are explained by underlying bipartite structures relating actors to movies and scientists to papers [7, 12]. Such

networks are naturally modeled using directed or undirected random intersection graphs [1,3,4,6,8].

A directed intersection graph on a node set  $V = \{1, \dots, n\}$  is constructed with the help of an auxiliary set of attributes  $W = \{w_1, \dots, w_m\}$  and a directed bipartite graph  $H$  with bipartition  $V \cup W$ , which models how nodes (or actors) relate to attributes. We say that actor  $i$  *demands* (or *follows*) attribute  $w_k$  when  $i \rightarrow w_k$ , and *supplies* it when  $i \leftarrow w_k$ . The directed intersection graph  $D$  induced by  $H$  is the directed graph on  $V$  such that  $i \rightarrow j$  if and only if  $H$  contains a path  $i \rightarrow w_k \rightarrow j$ , or equivalently,  $i$  demands one or more attributes supplied by  $j$  (see Fig. 2). For example, in a citation network the fact that an author  $i$  cites a paper  $w_k$  coauthored by  $j$ , corresponds to  $i \rightarrow w_k \rightarrow j$ .



**Fig. 2.** Node 1 follows node 2, because 1 demands attribute  $w_1$  supplied by 2.

We consider a random bipartite digraph  $H$  where the pairs  $(i, w_k)$ ,  $i \in V$ ,  $w_k \in W$  establish adjacency relations independently of each other. That is, the bivariate binary random vectors  $(\mathbb{I}_{i \rightarrow k}, \mathbf{I}_{k \rightarrow i})$ ,  $1 \leq i \leq n$ ,  $1 \leq k \leq m$ , are stochastically independent. Here  $\mathbb{I}_{i \rightarrow k}$  and  $\mathbf{I}_{k \rightarrow i}$  stand for the indicators of the events that links  $i \rightarrow w_k$  and  $w_k \rightarrow i$  are present in  $H$ . We assume that every pair  $(i, w_k)$  is assigned a triple of probabilities

$$p_{ik} = P(i \rightarrow w_k), \quad q_{ik} = P(w_k \rightarrow i), \quad r_{ik} = P(i \rightarrow w_k, w_k \rightarrow i). \quad (5)$$

Note that, by definition,  $r_{ik}$  satisfies the inequalities

$$\max\{p_{ik} + q_{ik} - 1, 0\} \leq r_{ik} \leq \min\{p_{ik}, q_{ik}\}. \quad (6)$$

A collection of triples  $\{(p_{ik}, q_{ik}, r_{ik}), 1 \leq i \leq n, 1 \leq k \leq m\}$  defines the distribution of a random bipartite digraph  $H$ .

We will focus on a fitness model where every node  $i$  is prescribed a pair of weights  $x_i, y_i \geq 0$  modelling the demand and supply intensities of  $i$ . Similarly, every attribute  $w_k$  is prescribed a weight  $z_k > 0$  modelling its relative popularity. Letting

$$p_{ik} = \min\{1, \gamma x_i z_k\} \quad \text{and} \quad q_{ik} = \min\{1, \gamma y_i z_k\}, \quad i, k \geq 1, \quad (7)$$

we obtain link probabilities proportional to respective weights. Furthermore, we assume that

$$r_{ik} = r(x_i, y_i, z_k, \gamma), \quad i, k \geq 1, \quad (8)$$

for some function  $r \geq 0$  satisfying (6). Here  $\gamma > 0$  is a parameter, defining the link density in  $H$ , which generally depends on  $m$  and  $n$ . Note that  $r$  defines the correlation between reciprocal links  $i \rightarrow w_k$  and  $w_k \rightarrow i$ . For example, by letting  $r(x, y, z, \gamma) = (\gamma xz \wedge 1)(\gamma yz \wedge 1)$ , we obtain a random bipartite digraph with independent links.

We will consider weight sequences having desired statistical properties for complex network modelling. For this purpose we assume that the node and attributes weights are realizations of random sequences  $X = (X_i)_{i \geq 1}$ ,  $Y = (Y_i)_{i \geq 1}$ , and  $Z = (Z_k)_{k \geq 1}$ , such that the sequences  $\{(X_i, Y_i), i \geq 1\}$  and  $\{Z_k, k \geq 1\}$  are mutually independent and consist of independent and identically distributed terms. The resulting random bipartite digraph is denoted by  $\mathcal{H}$ , and the resulting random intersection digraph by  $\mathcal{D}$ . We remark that  $\mathcal{D}$  extends the random intersection digraph model introduced in [1].

### 1.3 Degree distributions

When  $\gamma = (mn)^{-1/2}$  and  $m, n \rightarrow \infty$ , the random digraph  $\mathcal{D}$  defined in Sec. 1.2 becomes sparse, having the number of links proportional to the number of nodes. Theorem 1 below describes the class of limiting distributions of the outdegree of a typical vertex  $i$ . We remark that for each  $n$  the outdegrees  $d_+(1), \dots, d_+(n)$  are identically distributed.

To state the theorem, we let  $\Lambda_1, \Lambda_2, \Lambda_3$  be mixed-Poisson random variables distributed according to

$$P(\Lambda_i = r) = E e^{-\lambda_i} \frac{\lambda_i^r}{r!}, \quad r \geq 0,$$

where  $\lambda_1 = X_1 \beta^{1/2} E Z_1$ ,  $\lambda_2 = Z_1 \beta^{-1/2} E Y_1$ , and  $\lambda_3 = X_1 (E Y_1) (E Z_1^2)$ . We also denote by  $\Lambda_i^*$  a downshifted size-biased version of  $\Lambda_i$ , distributed according to

$$P(\Lambda_i^* = r) = \frac{r+1}{E \Lambda_i} P(\Lambda_i = r+1), \quad r \geq 0.$$

Below  $\xrightarrow{d}$  refers to convergence in distribution.

**Theorem 1.** *Consider a model with  $n, m \rightarrow \infty$  and  $\gamma = (nm)^{-1/2}$ , and assume that  $E Y_1, E Z_1^2 < \infty$ .*

- (i) *If  $m/n \rightarrow 0$  then  $d_+(1) \xrightarrow{d} 0$ .*
- (ii) *If  $m/n \rightarrow \beta$  for some  $\beta \in (0, \infty)$ , then  $d_+(1) \xrightarrow{d} \sum_{j=1}^{\Lambda_1} \Lambda_{2,j}^*$ , where  $\Lambda_{2,1}^*, \Lambda_{2,2}^*, \dots$  are independent copies of  $\Lambda_2^*$ , and independent of  $\Lambda_1$ .*
- (iii) *If  $m/n \rightarrow \infty$ , then  $d_+(1) \xrightarrow{d} \Lambda_3$ .*

*Remark 2.* By symmetry, the results of Theorem 1 extend to the indegree  $d_-(1)$  when we redefine  $\lambda_1 = Y_1 \beta^{1/2} E Z_1$ ,  $\lambda_2 = Z_1 \beta^{-1/2} E X_1$ , and  $\lambda_3 = Y_1 (E X_1) (E Z_1^2)$ .

*Remark 3.* The limiting distributions appearing in Theorem 1:(ii)–(iii) admit heavy tails. This random digraph model is rich enough to model power-law in-degree and outdegree distributions, or power-law indegree and light-tailed out-degree distributions.

*Remark 4.* The moment conditions in Theorem 1 are not the sharpest possible. For example, in (i) it is sufficient to assume that  $E Z_1 < \infty$ .

We note that a related result for simple (undirected) random intersection graph has been shown in [2]. Theorem 1 extends the result of [2] to digraphs.

#### 1.4 Diclique clustering

We investigate clustering in the random digraph  $\mathcal{D}$  defined in Sec. 1.2 by approximating the (random) diclique clustering coefficient  $\mathcal{C}_{\text{di}}(\mathcal{D})$  defined in (2) by a related nonrandom quantity

$$c_{\text{di}} := P(I_2 \rightarrow I_4 \mid I_1 \rightarrow I_3, I_1 \rightarrow I_4, I_2 \rightarrow I_3),$$

where  $(I_1, I_2, I_3, I_4)$  is a random ordered quadruple of distinct nodes chosen uniformly at random. Note that here  $P$  refers to two independent sources of randomness: the random digraph generation mechanism and the sampling of the nodes. Because the distribution of  $\mathcal{D}$  is invariant with respect to a relabeling of the nodes, the above quantity can also be written as

$$c_{\text{di}} = P(2 \rightarrow 4 \mid 1 \rightarrow 3, 1 \rightarrow 4, 2 \rightarrow 3).$$

We believe that under mild regularity conditions  $\mathcal{C}_{\text{di}}(\mathcal{D}) \approx c_{\text{di}}$ , provided that  $m$  and  $n$  are sufficiently large. Proving this is left for future work.

Theorem 2 below shows that the random digraph  $\mathcal{D}$  admits a *nonvanishing* clustering coefficient  $c_{\text{di}}$  when the intensity  $\gamma$  is inversely proportional to the number of attributes. For example, by choosing  $\gamma = (nm)^{-1/2}$  and letting  $m, n \rightarrow \infty$  so that  $m/n \rightarrow \beta > 0$ , we obtain a sparse random digraph with tunable clustering coefficient  $c_{\text{di}}$  and limiting degree distributions defined by Theorem 1 and Remark 2.

**Theorem 2.** Assume that  $m \rightarrow \infty$  and  $\gamma m \rightarrow \alpha$  for some constant  $\alpha \in (0, \infty)$ , and that  $E X_1^3, E Y_1^3, E Z_1^4 < \infty$ . Then

$$c_{\text{di}} \rightarrow \left( 1 + \alpha \left( \frac{E X_1^2}{E X_1} + \frac{E Y_1^2}{E Y_1} \right) \frac{(E Z_1^2)(E Z_1^3)}{E Z_1^4} + \alpha^2 \frac{E X_1^2}{E X_1} \frac{E Y_1^2}{E Y_1} \frac{(E Z_1^2)^3}{E Z_1^4} \right)^{-1}. \quad (9)$$

*Remark 5.* When  $E X_1^4, E Y_1^4, E Z_1^4 < \infty$ , the argument in the proof of Theorem 2 allows to conclude that  $c_{\text{di}} \rightarrow 0$  when  $\gamma m \rightarrow \infty$ .

To investigate clustering among the followers of a particular ego node  $i$ , we study a theoretical analogue of the local diclique clustering coefficient  $\mathcal{C}_{\text{di}}(D, i)$  defined in (4). By symmetry, we may relabel the nodes so that  $i = 3$ . We will consider the weights of node 3 as known and analyze the conditional probability

$$c_{\text{di}}(X_3, Y_3) = P_{X_3, Y_3}(2 \rightarrow 4 \mid 1 \rightarrow 3, 1 \rightarrow 4, 2 \rightarrow 3),$$

where  $P_{X_3, Y_3}$  refers to the conditional probability given  $(X_3, Y_3)$ . Actually, we may replace  $P_{X_3, Y_3}$  by  $P_{Y_3}$  above, because all events appearing on the right are independent of  $X_3$ .

One may also be interested in analyzing the conditional probability

$$c_{\text{di}}(X, Y) = P_{X, Y}(2 \rightarrow 4 \mid 1 \rightarrow 3, 1 \rightarrow 4, 2 \rightarrow 3),$$

where  $P_{X, Y}$  refers to the conditional probability given the values of all node weights  $X = (X_i)$  and  $Y = (Y_i)$ . Again, we may replace  $P_{X, Y}$  by  $P_{X_1, X_2, Y_3, Y_4}$  above, because the events on the right are independent of the other nodes' weights. More interestingly,  $c_{\text{di}}(X, Y)$  turns out to be asymptotically independent of  $X_2$  and  $Y_4$  as well in the sparse regime.

**Theorem 3.** *Assume that  $m \rightarrow \infty$  and  $\gamma m \rightarrow \alpha$  for some constant  $\alpha \in (0, \infty)$ .*

(i) *If  $E X_1^3, E Y_1^3, E Z_1^4 < \infty$ , then*

$$c_{\text{di}}(X_3, Y_3) \xrightarrow{P} \left( 1 + \alpha \left( \frac{E X_1^2}{E X_1} + Y_3 \right) \frac{(E Z_1^3)(E Z_1^2)}{E Z_1^4} + \alpha^2 Y_3 \frac{E X_1^2}{E X_1} \frac{(E Z_1^2)^3}{E Z_1^4} \right)^{-1}.$$

(ii) *If  $E Z_1^4 < \infty$ , then*

$$c_{\text{di}}(X, Y) \xrightarrow{P} \left( 1 + \alpha(X_1 + Y_3) \frac{(E Z_1^3)(E Z_1^2)}{E Z_1^4} + \alpha^2 X_1 Y_3 \frac{(E Z_1^2)^3}{E Z_1^4} \right)^{-1}.$$

Note that for large  $Y_3$ , the clustering coefficient  $c_{\text{di}}(X_3, Y_3) = c_{\text{di}}(Y_3)$  scales as  $Y_3^{-1}$ . Similarly, for large  $X_1$  and  $Y_3$ , the probability  $c_{\text{di}}(X, Y)$  scales as  $X_1^{-1} Y_3^{-1}$ . We remark that similar scaling of a related clustering coefficient in an undirected random intersection graph has been observed in [4].

*Remark 6.* When all attribute weights are equal to a constant  $z > 0$ , the statement in Theorem 3:(ii) simplifies into  $c_{\text{di}}(X, Y) \xrightarrow{P} (1 + \alpha z X_1)^{-1} (1 + \alpha z Y_3)^{-1}$ , a result reported in [9].

*Remark 7.* Theorems 1, 2, and 3 do not impose any restrictions on the correlation structure of the supply and demand indicators defined by (8).

### 1.5 Diclique versus transitivity clustering

An interesting question is to compare the diclique clustering coefficient  $c_{\text{di}}$  with the commonly used transitive closure clustering coefficient

$$c_{\text{tr}} = P(2 \rightarrow 4 \mid 2 \rightarrow 3 \rightarrow 4),$$

see e.g. [5, 15]. The next result illustrates that  $c_{\text{tr}}$  depends heavily on the correlation between the supply and demand indicators characterized by the function  $r(x, y, z, \gamma)$  in (8). A similar finding for a related random intersection graph has been discussed in [1]. We denote  $\min\{a, b\} = a \wedge b$ .

**Theorem 4.** *Let  $m, n \rightarrow \infty$ . Assume that  $\gamma = (nm)^{-1/2}$  and  $m/n \rightarrow \beta$  for some  $\beta > 0$ . Suppose also that  $E X_1^2, E Y_1^2, E Z_1^2 < \infty$ .*

- (i) *If  $r(x, y, z, \gamma) = (\gamma x z \wedge 1)(\gamma y z \wedge 1)$ , then  $c_{\text{tr}} \rightarrow 0$ .*
- (ii) *If  $r(x, y, z, \gamma) = \epsilon(\gamma x z \wedge \gamma y z \wedge 1)$  for some  $0 < \epsilon \leq 1$  and  $E(X_1 \wedge Y_1) > 0$ , then*

$$c_{\text{tr}} \rightarrow \left(1 + \frac{\sqrt{\beta}}{\epsilon} \frac{E(X_1 Y_1)}{E(X_1 \wedge Y_1)} \frac{(E Z_1^2)^2}{E Z_1^3}\right)^{-1}. \quad (10)$$

The assumption in (i) means that the supply and demand indicators of any particular node–attribute pair are conditionally independent given the weights. In contrast, the assumption in (ii) forces a strong correlation between the supply and demand indicators. We note that condition (6) is satisfied in case (ii) for all  $i \leq n$  and  $k \leq m$  with high probability as  $n, m \rightarrow \infty$ , because  $n^{-1/2} \max_{i \leq n} (X_i + Y_i) \xrightarrow{P} 0$  and  $m^{-1/2} \max_{k \leq m} Z_k \xrightarrow{P} 0$  imply that  $\gamma X_i Z_k + \gamma Y_i Z_k \leq 1$  for all  $i \leq n$  and  $k \leq m$  with high probability.

We remark that in case (i), and in case (ii) with a very small  $\epsilon$ , the transitive closure clustering coefficient  $c_{\text{tr}}$  becomes negligibly small, whereas the diclique clustering coefficient  $c_{\text{di}}$  remains bounded away from zero. Hence, it make sense to consider the event  $\{1 \rightarrow 3, 1 \rightarrow 4, 2 \rightarrow 3\}$  as a more robust predictor of the link  $2 \rightarrow 4$  than the event  $\{2 \rightarrow 3 \rightarrow 4\}$ . This conclusion has been empirically confirmed for various real-world networks in [10, 17].

## 2 Proofs

The proof of Theorem 1 goes along similar lines as that of Theorem 1 in [2]. It is omitted. We only give the proofs of Theorems 2 and 3. The proof of Theorem 4 is given in an extended version of the paper available from the authors.

We assume for notational convenience that  $\gamma = \alpha m^{-1}$ . Denote events  $\mathcal{A} = \{1 \rightarrow 3, 1 \rightarrow 4, 2 \rightarrow 3\}$ ,  $\mathcal{B} = \{2 \rightarrow 4\}$  and random variables

$$\tilde{p}_{ik} = \alpha \frac{X_i Z_k}{m}, \quad \tilde{q}_{ik} = \alpha \frac{Y_i Z_k}{m}.$$

By  $\tilde{P}$  and  $\tilde{E}$  we denote the conditional probability and expectation given  $X, Y, Z$ . Note that  $p_{ik} = \tilde{P}(\mathbb{I}_{i \rightarrow k} = 1)$ ,  $q_{ik} = \tilde{P}(\mathbf{I}_{k \rightarrow i} = 1)$ , and

$$p_{ik} = 1 \wedge \tilde{p}_{ik}, \quad q_{ik} = 1 \wedge \tilde{q}_{ik}. \quad (11)$$

*Proof (of Theorem 2).* We observe that  $\mathcal{A} = \cup_{i \in [4]} \mathcal{A}_i$ , where

$$\begin{aligned}\mathcal{A}_1 &= \bigcup_{k \in \mathcal{C}_1} \mathcal{A}_{1.k}, & \mathcal{A}_{1.k} &= \{\mathbb{I}_{1 \rightarrow k} \mathbb{I}_{2 \rightarrow k} \mathbf{I}_{k \rightarrow 3} \mathbf{I}_{k \rightarrow 4} = 1\}, \\ \mathcal{A}_2 &= \bigcup_{(k,l) \in \mathcal{C}_2} \mathcal{A}_{2.kl}, & \mathcal{A}_{2.kl} &= \{\mathbb{I}_{1 \rightarrow k} \mathbb{I}_{2 \rightarrow l} \mathbf{I}_{k \rightarrow 3} \mathbf{I}_{k \rightarrow 4} \mathbf{I}_{l \rightarrow 3} = 1\}, \\ \mathcal{A}_3 &= \bigcup_{(k,l) \in \mathcal{C}_3} \mathcal{A}_{3.kl}, & \mathcal{A}_{3.kl} &= \{\mathbb{I}_{1 \rightarrow k} \mathbb{I}_{1 \rightarrow l} \mathbb{I}_{2 \rightarrow k} \mathbf{I}_{k \rightarrow 3} \mathbf{I}_{l \rightarrow 4} = 1\}, \\ \mathcal{A}_4 &= \bigcup_{(j,k,l) \in \mathcal{C}_4} \mathcal{A}_{4.jkl}, & \mathcal{A}_{4.jkl} &= \{\mathbb{I}_{1 \rightarrow j} \mathbb{I}_{1 \rightarrow k} \mathbb{I}_{2 \rightarrow l} \mathbf{I}_{j \rightarrow 3} \mathbf{I}_{k \rightarrow 4} \mathbf{I}_{l \rightarrow 3} = 1\}.\end{aligned}$$

Here  $\mathcal{C}_1 = [m]$ ,  $\mathcal{C}_2 = \mathcal{C}_3 = \{(k, l) : k \neq l; k, l \in [m]\}$ , and  $\mathcal{C}_4 = \{(j, k, l) : j \neq k \neq l; j, k, l \in [m]\}$ . Hence, by inclusion-exclusion,

$$\sum_{i \in [4]} P(\mathcal{A}_i) - \sum_{\{i, j\} \subset [4]} P(\mathcal{A}_i \cap \mathcal{A}_j) \leq P(\mathcal{A}) \leq \sum_{i \in [4]} P(\mathcal{A}_i).$$

We prove the theorem in Claims 1 – 3 below. Claim 2 implies that  $P(\mathcal{A}) = \sum_{i \in [4]} P(\mathcal{A}_i) + O(m^{-4})$ . Claim 3 implies that  $P(\mathcal{A} \cap \mathcal{B}) = P(\mathcal{A}_1) + O(m^{-4})$ . Finally, Claim 1 establishes the approximation (9) to the ratio  $\mathcal{C}_{\text{di}} = P(\mathcal{A} \cap \mathcal{B})/P(\mathcal{A})$ .

*Claim 1. We have*

$$P(\mathcal{A}_1) = \alpha^4 m^{-3} A_1 (1 + o(1)), \quad (12)$$

$$P(\mathcal{A}_2) = \alpha^5 m^{-3} A_2 (1 + o(1)), \quad (13)$$

$$P(\mathcal{A}_3) = \alpha^5 m^{-3} A_3 (1 + o(1)), \quad (14)$$

$$P(\mathcal{A}_4) = \alpha^6 m^{-3} A_4 (1 + o(1)). \quad (15)$$

Here we denote

$$A_1 = a_1^2 b_1^2 h_4, \quad A_2 = a_1^2 b_1 b_2 h_2 h_3, \quad A_3 = a_1 a_2 b_1^2 h_2 h_3, \quad A_4 = a_1 a_2 b_1 b_2 h_2^3.$$

and  $a_r = E X_1^r$ ,  $b_4 = E Y_1^r$ ,  $h_r = E Z_1^r$ .

*Claim 2.* For  $1 \leq i < j \leq 4$  we have

$$P(\mathcal{A}_i \cap \mathcal{A}_j) = O(m^{-4}). \quad (16)$$

*Claim 3.* We have

$$P(\mathcal{B} \cap \mathcal{A}) = P(\mathcal{A}_1) + O(m^{-4}). \quad (17)$$

*Proof of Claim 1.* We estimate every  $P(\mathcal{A}_r)$  using inclusion-exclusion  $I_1 - I_2 \leq P(\mathcal{A}_r) \leq I_1$ . Here

$$I_1 = I_1(r) = \sum_{x \in \mathcal{C}_r} P(\mathcal{A}_{r.x}), \quad I_2 = I_2(r) = \sum_{\{x, y\} \subset \mathcal{C}_r} P(\mathcal{A}_{r.x} \cap \mathcal{A}_{r.y}).$$



Now (12-15) follow from the approximations

$$\begin{aligned} I_1 &= \alpha^4 m^{-3} A_1(1 + o(1)), & I_2 &= \alpha^5 m^{-3} A_2(1 + o(1)), \\ I_3 &= \alpha^5 m^{-3} A_3(1 + o(1)), & I_4 &= \alpha^6 m^{-3} A_4(1 + o(1)) \end{aligned} \quad (18)$$

and bounds  $I_2(r) = o(m^{-3})$ , for  $1 \leq r \leq 4$ .

Firstly we show (18). We only prove the first relation. The remaining cases are treated in much the same way. From the inequalities, see (11),

$$\begin{aligned} \tilde{p}_{1k} \tilde{p}_{2k} \tilde{q}_{3k} \tilde{q}_{4k} &\geq p_{1k} p_{2k} q_{3k} q_{4k} \geq \tilde{p}_{1k} \tilde{p}_{2k} \tilde{q}_{3k} \tilde{q}_{4k} \mathbb{I}'_k \geq \tilde{p}_{1k} \tilde{p}_{2k} \tilde{q}_{3k} \tilde{q}_{4k} - \tilde{p}_{1k} \tilde{p}_{2k} \tilde{q}_{3k} \tilde{q}_{4k} \mathbb{I}^*_k, \\ \mathbb{I}'_k &= \mathbb{I}_{\tilde{p}_{1k} \leq 1} \mathbb{I}_{\tilde{p}_{2k} \leq 1} \mathbb{I}_{\tilde{q}_{3k} \leq 1} \mathbb{I}_{\tilde{q}_{4k} \leq 1}, \quad \mathbb{I}^*_k = \mathbb{I}_{\tilde{p}_{1k} > 1} + \mathbb{I}_{\tilde{p}_{2k} > 1} + \mathbb{I}_{\tilde{q}_{3k} > 1} + \mathbb{I}_{\tilde{q}_{4k} > 1}, \end{aligned}$$

we obtain that

$$P(\mathcal{A}_{1.k}) = E p_{1k} p_{2k} q_{3k} q_{4k} = E \tilde{p}_{1k} \tilde{p}_{2k} \tilde{q}_{3k} \tilde{q}_{4k} + R, \quad (19)$$

where

$$E \tilde{p}_{1k} \tilde{p}_{2k} \tilde{q}_{3k} \tilde{q}_{4k} = \alpha^4 m^{-4} A_1 \quad \text{and} \quad |R| \leq E \tilde{p}_{1k} \tilde{p}_{2k} \tilde{q}_{3k} \tilde{q}_{4k} \mathbb{I}^*_k = o(m^{-4}).$$

Hence  $I_1 = mP(\mathcal{A}_{1.k}) = \alpha^4 m^{-3} A_1(1 + o(1))$ .

Secondly we show that  $I_2(r) = o(m^{-3})$ , for  $1 \leq r \leq 4$ . For  $r = 1$  the bound  $I_2(1) = \binom{m}{2} P(\mathcal{A}_{1.k} \cap \mathcal{A}_{1.l}) = o(m^{-3})$  follows from the inequalities

$$P(\mathcal{A}_{1.k} \cap \mathcal{A}_{1.l}) \leq E \tilde{p}_{1k} \tilde{p}_{2k} \tilde{q}_{3k} \tilde{q}_{4k} \tilde{p}_{1l} \tilde{p}_{2l} \tilde{q}_{3l} \tilde{q}_{4l} = O(m^{-8}).$$

For  $r = 2, 3$  we split  $I_2(r) = J_1 + \dots + J_5$ , where

$$\begin{aligned} J_1 &= \sum_{\{(k,l),(k',l')\} \subset \mathcal{C}_r} P(\mathcal{A}_{r.kl} \cap \mathcal{A}_{r.k'l'}), & J_2 &= \sum_{\{(k,l),(k',l')\} \subset \mathcal{C}_r} P(\mathcal{A}_{r.kl} \cap \mathcal{A}_{r.k'l}), \\ J_3 &= \sum_{\{(k,l),(k',l')\} \subset \mathcal{C}_r} P(\mathcal{A}_{r.kl} \cap \mathcal{A}_{r.k'l'}), & J_4 &= \sum_{\{(k,l),(k',k)\} \subset \mathcal{C}_r, k' \neq l} P(\mathcal{A}_{r.kl} \cap \mathcal{A}_{r.k'k}), \\ J_5 &= \sum_{(k,l) \in \mathcal{C}_r} P(\mathcal{A}_{r.kl} \cap \mathcal{A}_{r.lk}). \end{aligned}$$

In the first (second) sum distinct pairs  $x = (k, l)$  and  $y = (k', l')$  share the first (second) coordinate. In the third sum all coordinates of the pairs  $(k, l), (k', l')$  are different. In the fourth sum the pairs  $(k, l), (k', k)$  only share one common element, but it appears in different coordinates. We show that each  $J_i = o(m^{-3})$ . Next we only consider the case of  $r = 2$ . The case of  $r = 3$  is treated in a similar

way. We have

$$\begin{aligned}
J_1 &= m \binom{m-1}{2} P(\mathcal{A}_{2.kl} \cap \mathcal{A}_{2.kl'}) \leq m^3 E H_1, \quad H_1 = p_{1k} p_{2l} p_{2l'} q_{3k} q_{4k} q_{3l} q_{3l'}, \\
J_2 &= m \binom{m-1}{2} P(\mathcal{A}_{2.kl} \cap \mathcal{A}_{2.k'l}) \leq m^3 E H_2, \quad H_2 = p_{1k} p_{1k'} p_{2l} q_{3k} q_{4k} q_{3k'} q_{4k'} q_{3l}, \\
J_3 &= \binom{m}{2} \binom{m-2}{2} P(\mathcal{A}_{2.kl} \cap \mathcal{A}_{2.k'l'}) \leq m^4 E H_3, \quad H_3 = p_{1k} p_{1k'} p_{2l} p_{2l'} q_{3k} q_{4k} q_{3k'} q_{4k'} q_{3l} q_{3l'}, \\
J_4 &= m(m-1)(m-2) P(\mathcal{A}_{2.kl} \cap \mathcal{A}_{2.k'k}) \leq m^3 E H_4, \quad H_4 = p_{1k} p_{1k'} p_{2l} p_{2k} q_{3k} q_{4k} q_{3k'} q_{4k'} q_{3l}, \\
J_5 &= \binom{m}{2} P(\mathcal{A}_{2.kl} \cap \mathcal{A}_{2.lk}) \leq m^2 E H_5, \quad H_5 = p_{1k} p_{1l} p_{2l} p_{2k} q_{3k} q_{3l} q_{4k} q_{4l}.
\end{aligned}$$

In the product  $H_1$  we estimate the typical factors  $p_{ij} \leq \tilde{p}_{ij}$  and  $q_{ij} \leq \tilde{q}_{ij}$ , but

$$q_{3l} \leq \tilde{q}_{3l} \mathbb{I}_{Y_3 \leq \sqrt{m}} + \mathbb{I}_{Y_3 > \sqrt{m}} \leq \alpha m^{-1/2} Z_l + \mathbb{I}_{Y_3 > \sqrt{m}}. \quad (20)$$

We obtain

$$E H_1 \leq \alpha^6 m^{-6} a_1 a_2 b_1 h_2 h_3 (b_2 h_2 \alpha m^{-1/2} + h_1 E Y_3^2 \mathbb{I}_{Y_3 > \sqrt{m}}) = o(m^{-6}). \quad (21)$$

Hence  $J_1 = o(m^{-3})$ . Similarly, we show that  $J_2 = o(m^{-4})$ . Furthermore, while estimating  $H_3$  we apply (20) to  $q_{3l}$  and  $q_{3l'}$  and apply  $p_{ij} \leq \tilde{p}_{ij}$  and  $q_{ij} \leq \tilde{q}_{ij}$  to remaining factors. We obtain

$$H_3 \leq \tilde{p}_{1k} \tilde{p}_{1k'} \tilde{p}_{2l} \tilde{p}_{2l'} \tilde{q}_{3k} \tilde{q}_{4k} \tilde{q}_{3k'} \tilde{q}_{4k'} (\alpha m^{-1/2} Z_l + \mathbb{I}_{Y_3 > \sqrt{m}}) (\alpha m^{-1/2} Z_{l'} + \mathbb{I}_{Y_3 > \sqrt{m}}). \quad (22)$$

Since the expected value of the product on the right is  $o(m^{-8})$ , we conclude that  $E H_3 = o(m^{-8})$ . Hence  $J_3 = o(m^{-4})$ . Proceeding in a similar way we establish the bounds  $J_4 = o(m^{-5})$  and  $J_5 = O(m^{-6})$ .

We explain the truncation step (20) in some more detail. A simple upper bound for  $H_1$  is the product

$$\tilde{p}_{1k} \tilde{p}_{2l} \tilde{p}_{2l'} \tilde{q}_{3k} \tilde{q}_{4k} \tilde{q}_{3l} \tilde{q}_{3l'} = \alpha^7 m^{-7} X_1 X_2^2 Y_3^3 Y_4 Z_k^3 Z_l^2 Z_{l'}^2.$$

It contains an undesirable high power  $Y_3^3$ . Using (20) instead of the simple upper bound  $q_{3l} \leq \tilde{q}_{3l}$  we have reduced in (21) the power of  $Y_3$  down to 2. Similarly, in (22) we have reduced the power of  $Y_3$  from 4 to 2.

Using the truncation argument we obtain the upper bound  $I_2(4) = o(m^{-3})$  under moment conditions  $E X_1^3, E Y_1^3, E Z_1^4 < \infty$ . The proof is similar to that of the bound  $I_2(2) = o(m^{-3})$  above. We omit routine, but tedious calculation.

*Proof of Claim 2.* We only prove that  $q := P(\mathcal{A}_3 \cap \mathcal{A}_4) = O(m^{-4})$ . The remaining cases are treated in a similar way. For  $x = (j, k, l) \in \mathcal{C}_4$  and  $y = (r, t) \in \mathcal{C}_3$  we denote, for short,  $\mathbb{I}_{\mathcal{A}_{4,x}} = \mathbb{I}_x^* = \mathbb{I}_{jkl}^*$  and  $\mathbb{I}_{\mathcal{A}_{3,y}} = \mathbb{I}_y = \mathbb{I}_{rt}$ . For  $q = E \mathbb{I}_{\mathcal{A}_4} \mathbb{I}_{\mathcal{A}_3}$ , we write, by the symmetry,

$$q \leq E \left( \sum_{x \in \mathcal{C}_4} \mathbb{I}_x^* \right) \mathbb{I}_{\mathcal{A}_3} = m(m-1)(m-2) E \mathbb{I}_{123}^* \mathbb{I}_{\mathcal{A}_3}$$

and

$$E \mathbb{I}_{123}^* \mathbb{I}_{\mathcal{A}_3} \leq E \mathbb{I}_{123}^* \left( \sum_{y \in \mathcal{C}_3} \mathbb{I}_y \right) = E \mathbb{I}_{123}^* (J_1 + J_2 + J_3).$$

Here

$$J_1 = \sum_{r,t \in [m] \setminus [3], r \neq t} \mathbb{I}_{rt}, \quad J_2 = \sum_{r \in [m] \setminus [3]} \sum_{s \in [3]} (\mathbb{I}_{sr} + \mathbb{I}_{rs}), \quad J_3 = \sum_{r,t \in [3], r \neq t} \mathbb{I}_{rt}.$$

Finally, we show that  $E \mathbb{I}_{123}^* J_i = O(m^{-7})$ ,  $i \in [3]$ . For  $i = 1$  we have, by the symmetry,

$$E \mathbb{I}_{123}^* J_1 = (m-3)(m-4)E \mathbb{I}_{123}^* \mathbb{I}_{45}. \quad (23)$$

Invoking the inequalities

$$E \mathbb{I}_{123}^* \mathbb{I}_{45} = E \tilde{E} \mathbb{I}_{123}^* \mathbb{I}_{45} \leq E \tilde{p}_{11} \tilde{p}_{12} \tilde{p}_{15} \tilde{p}_{23} \tilde{p}_{24} \tilde{q}_{13} \tilde{q}_{24} \tilde{q}_{33} \tilde{q}_{43} \tilde{q}_{54} = O(m^{-10}) \quad (24)$$

we obtain  $E \mathbb{I}_{123}^* J_1 = O(m^{-8})$ .

The bound  $E \mathbb{I}_{123}^* J_2 = O(m^{-7})$  is obtained from the identity (which follows by symmetry)

$$E \mathbb{I}_{123}^* J_2 = (m-3) \sum_{s \in [3]} (E \mathbb{I}_{123}^* \mathbb{I}_{s4} + E \mathbb{I}_{123}^* \mathbb{I}_{4s}),$$

combined with bounds  $E \mathbb{I}_{123}^* \mathbb{I}_{s4} + E \mathbb{I}_{123}^* \mathbb{I}_{4s} = O(m^{-8})$ ,  $s \in [3]$ . We only show the latter bound for  $s = 3$ . The cases  $s = 1, 2$  are treated in a similar way. We have

$$\begin{aligned} E \mathbb{I}_{123}^* \mathbb{I}_{34} &\leq E \tilde{p}_{11} \tilde{p}_{12} \tilde{p}_{13} \tilde{p}_{23} \tilde{q}_{13} \tilde{q}_{24} \tilde{q}_{33} \tilde{q}_{44} = O(m^{-8}), \\ E \mathbb{I}_{123}^* \mathbb{I}_{43} &\leq E \tilde{p}_{11} \tilde{p}_{12} \tilde{p}_{13} \tilde{p}_{23} \tilde{p}_{24} \tilde{q}_{13} \tilde{q}_{24} \tilde{q}_{33} \tilde{q}_{34} \tilde{q}_{43} = O(m^{-10}). \end{aligned}$$

The proof of  $E \mathbb{I}_{123}^* J_3 = O(m^{-7})$  is similar. It is omitted.

*Proof of Claim 3.* We use the notation  $\bar{\mathbb{I}}_{\mathcal{A}_j} = 1 - \mathbb{I}_{\mathcal{A}_j}$  for the indicator of the event  $\bar{\mathcal{A}}_j$  complement to  $\mathcal{A}_j$ . For  $2 \leq i \leq 4$  we denote  $\mathcal{H}_i = (\mathcal{A}_i \cap \mathcal{B}) \setminus \cup_{1 \leq j \leq i-1} \mathcal{A}_j$ . We have

$$P(\mathcal{A} \cap \mathcal{B}) = P(\cup_{i \in [4]} \mathcal{A}_i \cap \mathcal{B}) = P(\mathcal{A}_1 \cap \mathcal{B}) + R, \quad 0 \leq R \leq P(\cup_{2 \leq i \leq 4} \mathcal{H}_i).$$

Note that  $P(\mathcal{A}_1 \cap \mathcal{B}) = P(\mathcal{A}_1)$ . It remains to show that  $P(\mathcal{H}_i) = O(m^{-4})$ ,  $2 \leq i \leq 4$ .

We have, by the symmetry,

$$P(\mathcal{H}_2) = E \mathbb{I}_{\mathcal{A}_2} \mathbb{I}_{\mathcal{B}} \bar{\mathbb{I}}_{\mathcal{A}_1} \leq E \sum_{x \in \mathcal{C}_2} \mathbb{I}_{\mathcal{A}_{2,x}} \mathbb{I}_{\mathcal{B}} \bar{\mathbb{I}}_{\mathcal{A}_1} = m(m-1)E \mathbb{I}_{\mathcal{A}_{2,12}} \mathbb{I}_{\mathcal{B}} \bar{\mathbb{I}}_{\mathcal{A}_1}. \quad (25)$$

Furthermore, we have  $\mathbb{I}_{\mathcal{A}_{2,12}} \mathbb{I}_{\mathcal{B}} \bar{\mathbb{I}}_{\mathcal{A}_1} \leq \mathbb{I}_{\mathcal{A}_{2,12}} (\mathbf{I}_{2 \rightarrow 4} + \sum_{3 \leq j \leq m} \mathbb{I}_{2 \rightarrow j} \mathbf{I}_{j \rightarrow 4})$  and, by the symmetry,

$$E \mathbb{I}_{\mathcal{A}_{2,12}} \mathbb{I}_{\mathcal{B}} \bar{\mathbb{I}}_{\mathcal{A}_1} \leq E \mathbb{I}_{\mathcal{A}_{2,12}} \mathbf{I}_{2 \rightarrow 4} + (m-2)E \mathbb{I}_{\mathcal{A}_{2,12}} \mathbb{I}_{2 \rightarrow 3} \mathbf{I}_{3 \rightarrow 4}.$$

A simple calculation shows that  $E \mathbb{I}_{\mathcal{A}_{2,12}} \mathbf{I}_{2 \rightarrow 4} \leq E \tilde{p}_{11} \tilde{p}_{22} \tilde{q}_{13} \tilde{q}_{14} \tilde{q}_{23} \tilde{q}_{24} = O(m^{-6})$ . Similarly,  $E \mathbb{I}_{\mathcal{A}_{2,12}} \mathbb{I}_{2 \rightarrow 3} \mathbf{I}_{3 \rightarrow 4} = O(m^{-7})$ . Therefore,  $E \mathbb{I}_{\mathcal{A}_{2,12}} \mathbb{I}_{\mathcal{B}} \mathbb{I}_{\mathcal{A}_1} = O(m^{-6})$ . Now (25) implies  $P(\mathcal{H}_2) = O(m^{-4})$ . The bounds  $P(\mathcal{H}_j) = O(m^{-4})$ ,  $j = 3, 4$  are obtained in a similar way.

*Proof (of Theorem 3).* The proof is the same as that of Theorem 2, but while evaluating the probabilities of events  $\mathcal{A}$  and  $\mathcal{A} \cap \mathcal{B}$  we treat  $X_1, X_2, Y_3, Y_4$ , respectively  $Y_3$ , as constants.

## References

1. Bloznelis, M.: A random intersection digraph: Indegree and outdegree distributions. Discrete Math. 310(19), 2560–2566 (2010), <http://dx.doi.org/10.1016/j.disc.2010.06.018>
2. Bloznelis, M., Damarackas, J.: Degree distribution of an inhomogeneous random intersection graph. Electron. J. Combin. 20(3) (2013)
3. Bloznelis, M., Godehardt, E., Jaworski, J., Kurauskas, V., Rybarczyk, K.: Recent Progress in Complex Network Analysis: Properties of Random Intersection Graphs, pp. 79–88. Springer, Berlin, Heidelberg (2015), [http://dx.doi.org/10.1007/978-3-662-44983-7\\_7](http://dx.doi.org/10.1007/978-3-662-44983-7_7)
4. Deijfen, M., Kets, W.: Random intersection graphs with tunable degree distribution and clustering. Probab. Eng. Inform. Sc. 23(4), 661–674 (2009), <http://dx.doi.org/10.1017/S0269964809990064>
5. Fagiolo, G.: Clustering in complex directed networks. Phys. Rev. E 76, 026107 (Aug 2007), <http://link.aps.org/doi/10.1103/PhysRevE.76.026107>
6. Frieze, A., Karoński, M.: Introduction to Random Graphs. Cambridge University Press (2016)
7. Guillaume, J.L., Latapy, M.: Bipartite structure of all complex networks. Information Processing Letters 90(5), 215–221 (2004)
8. Karoński, M., Scheinerman, E.R., Singer-Cohen, K.B.: On random intersection graphs: The subgraph problem. Combin. Probab. Comput. 8(1-2), 131–159 (1999), <http://dx.doi.org/10.1017/S0963548398003459>
9. Leskelä, L.: Directed random intersection graphs. Presentation at 18th INFORMS Applied Probability Society Conference, Istanbul, Turkey (July 2015)
10. Milo, R., Shen-Orr, S., Itzkovitz, S., Kashtan, N., Chklovskii, D., Alon, U.: Network motifs: Simple building blocks of complex networks. Science 298(5594), 824–827 (2002), <http://science.sciencemag.org/content/298/5594/824>
11. Newman, M.E.J.: The structure and function of complex networks. SIAM Review 45(2), 167–256 (2003), <http://dx.doi.org/10.1137/S003614450342480>
12. Newman, M.E.J., Strogatz, S.H., Watts, D.J.: Random graphs with arbitrary degree distributions and their applications. Phys. Rev. E 64 (2001)
13. Scott, J.: Social Network Analysis. SAGE Publications (2012)
14. Szabó, G., Alava, M., Kertész, J.: Clustering in Complex Networks, pp. 139–162. Springer (2004), [http://dx.doi.org/10.1007/978-3-540-44485-5\\_7](http://dx.doi.org/10.1007/978-3-540-44485-5_7)
15. Wasserman, S., Faust, K.: Social Network Analysis: Methods and Applications. Cambridge University Press (1994)
16. Watts, D.J., Strogatz, S.H.: Collective dynamics of 'small-world' networks. Nature 393, 440–442 (1998)
17. Zhang, Q.M., Lü, L., Wang, W.Q., Zhu, Y.X., Tao, Z.: Potential theory for directed networks. PLoS ONE 8(2) (2013)